

GPU, and multiple GPUs based on the actual data measured from fNIRS.

複数の GPU を用いた 3次元 Smith-Waterman 法の性能評価

須戸里織^{†1} 吉見真聡^{†1}
三木光範^{†1} 廣安知之^{†1}

近年、画像処理専用のハードウェアである GPU (Graphics Processing Unit) を汎用計算に用いる動きが多くみられるようになってきている。GPU のような並列計算ハードウェアでの高速化が期待されているアプリケーションの中のひとつに、SW (Smith-Waterman) 法がある。SW 法は DNA 配列の相同性検索に用いられているが、応用研究として時系列データの解析にも用いられており、そのひとつに、うつ病診断に用いられる fNIRS による脳血流量の解析がある。fNIRS は多数の計測点から出力があり、2つの文字列の比較を行う SW 法を利用して解析するためには、一度に比較できる系列数の拡張が課題となる。そこで本研究報告では、3次元 SW 法を提案し、fNIRS のデータに基づいて、処理能力を CPU と GPU で、また GPU の枚数を変化させ評価した。

A performance evaluation of three-dimensional Smith-Waterman method with multiple GPUs

SAORI SUDO,^{†1} MASATO YOSHIMI,^{†1} MITSUNORI MIKI^{†1}
and TOMOYUKI HIROYASU^{†1}

Graphic processing units (GPUs), which was entirely used in image processing, have been applied to a lot of studies and products in general purpose computation. One of the most famous algorithm expected to accelerate by utilizing GPU is Smith-Waterman (SW) algorithm in bioinformatics. The algorithm is also applied to analysis of time-series data, such as cerebral blood flow detected by fNIRS which is used to make a diagnosis for depression. As fNIRS outputs data from a lot of measure points, hundreds of executions of SW-algorithm are required to obtain relationships among channels. The expansion of the number of series which analyzed simultaneously may be a solution to increase the analysis quality. This report proposed the three-dimensional SW algorithm as a first step to address the problem, and evaluated performances on a CPU, a

1. はじめに

近年、画像処理用のハードウェアである GPU (Graphic Processing Unit) は、本来の利用目的だけでなく汎用計算に利用されるようになってきた。GPU を画像処理以外の汎用計算に応用することを GPGPU (General-purpose computing on GPU)¹⁾ という。GPU は CPU と比べて、価格と消費電力量あたりの処理能力が高いことから注目を集めている。GPU は多くのコア数をもつため、条件分岐のない単純な計算を並列処理するのに性能を発揮するが、条件分岐を行う複雑な命令では処理能力が低下するという特徴がある。

SW (Smith-Waterman) 法²⁾ は、生物学の分野において DNA 配列中で同じ順序で並んでいる配列を見つける配列アラインメント³⁾ の手法のひとつである。SW 法は計算において細粒度並列性が高く、GPU や FPGA などの並列計算ハードウェアで高速化が行われてきた。SW 法は DNA 配列の解析のほか、うつ病の発見に効果を発揮する fNIRS⁴⁾ のデータ解析への利用が考えられている⁵⁾。SW 法は部分的な類似部分を探索するのに適しているため、fNIRS の解析に向いているとされる。SW 法は2つの文字列の比較を行うので、多数の時系列データが出力として得られる fNIRS の解析を行うには、全組み合わせについて SW 法を実行しなければならない。効率的に SW 法を用いて fNIRS のデータ解析を行うためには、SW 法で同時に比較できる文字データの個数を増やす必要がある。また、SW 法は比較する文字列が長くなるほど計算に時間がかかるため、GPU や FPGA などで高速化を行う必要がある。

そこで本研究報告では、比較できるデータ数を3つに拡張した3次元 SW 法の提案を行い、GPGPU 向けの統合開発環境である CUDA を用いて、GPU で高速に計算する手法について述べる。また、提案したアルゴリズムを CPU と GPU で、また、GPU の枚数を変化させて比較検討を行う。

以降では、2章と3章でそれぞれ SW 法と3次元 SW 法について、4章で CUDA の概要について、5章で関連研究について述べる。次に、6章で本研究報告における実装を述べ

^{†1} 同志社大学
Doshisha University

X \ Y	-	A	C	A	C
-	0	0	0	0	0
A	0				
G	0				
C	0				
A	0				

図1 文字列テーブル

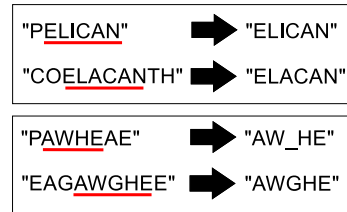


図2 SW法で抽出した類似部分列の例

て、7章で処理性能について評価を行い、8章でまとめを述べる。

2. Smith-Waterman 法

2.1 概要

SW法は2つの文字列の共通部分列を抽出するための手法である。SW法は部分的な類似部分を探る、局所的な配列アラインメントのひとつである。SW法は他の配列アラインメントと違い、類似度が最大の共通部分列のみでなく、複数の類似部分列を抽出することができる。

2.1.1 アルゴリズム

図1にSW法の計算で使用する文字列テーブルを示す。SW法では類似部分列を求めるために、図1に示すような対象文字列X, Yを、行と列に割り当てた文字列テーブルを計算する。図1の左上の要素から右下の要素へ、注目している2つの文字がどの程度類似しているかを示すスコアを計算し値を格納していく。文字列テーブルの右下まで計算を終えたのち、最も高いスコアの要素からスコアが0になるまで、要素を選択した経路を辿り、類似文字列を抽出する。この操作をトレースバックという。SW法によって抽出された類似文字列の例を図2に示す。

2.2 計算手順

SW法のスコアの計算ではmatch, mismatch, gapというパラメータを用いる。matchは2つの文字列が一致したとき、mismatchは不一致のとき、gapはスペースの挿入に関するスコアである。SW法では一般的にmatch = 1, mismatch = -1, gap = -1が用いられる。DNAの解析においては、BLOSUM行列と呼ばれる、2つのアミノ酸間の相同性スコアを示すアミノ酸置換テーブルの値を用いるが⁶⁾、本稿では時系列データの類似度を調べる

ため、アミノ酸置換テーブルは用いない。

長さがmの配列 $X = x_1x_2...x_m$ および長さがnである配列 $Y = y_1y_2...y_n$ を対象としたときの、SW法のアルゴリズムの流れを以下に示す。

Step1. 配列X, Yについて、文字列テーブルを作成し、それぞれの文字列を列と行に割り当てる。

Step2. i行j列にある要素がSW(i, j)であるとする、 $0 \leq k \leq m$ および $0 \leq l \leq n$ において、 $SW(k, l) = SW(k, 0) = SW(0, l) = 0$ と初期化を行う。

Step3. それぞれの要素について式(1)を基にスコアを計算する。

$$SW(i, j) = \max \begin{cases} SW(i-1, j-1) + A \\ SW(i-1, j) + \text{gap} \\ SW(i, j-1) + \text{gap} \\ 0 \end{cases} \quad (1)$$

$$A : \begin{cases} \text{match} & \text{if}(x_i = y_j) \\ \text{mismatch} & \text{else} \end{cases}$$

Step4. 最も高いスコアを持つ要素からトレースバックを行う。

文字列“ACAC”と“AGCA”の類似部分をSW法で求めたときの文字列テーブルを図3に示す。類似部分を抽出するために、最大のスコアの要素からスコアが0の要素までトレースバックを行う。

トレースバックの操作を図4に示す。図4において、最大のスコアを持つ要素は $(i, j) = (4, 3)$ となり、トレースバックは $(i, j) = (4, 3), (i, j) = (3, 2), (i, j) = (2, 1)$ という経路をたどる。スコアが0になるとトレースバックは終了する。トレースバックでは、それぞれの要素がスコアを計算するさいに使用した要素を逆向きに辿り、スコアが0のところは含まずに、要素の座標に対応している文字列から類似部分列を抽出する。“ACAC”と“AGCA”からは“CA”が類似度の高い部分文字列として抽出される。

3. 3次元 Smith-Waterman 法

3.1 概要

3つ以上の文字列の類似部分列を求める手法には様々なものがあるが、そのほとんどが近似法を用いており、複数の類似部分列を抽出することができない。複数の類似部分列を取得するにはSW法が適しているが、2つの文字列の比較が限度である。本研究報告では、SW法で3つの文字列の比較を行えるように、3次元SW法を作成する。3次元SW法で

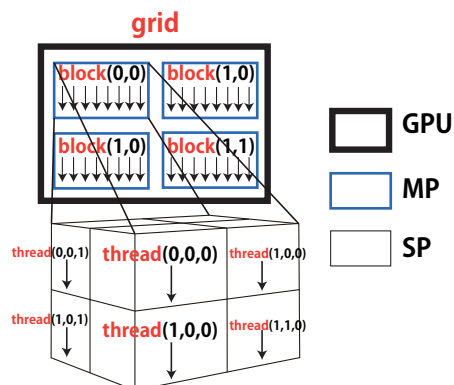


図 6 CUDA における GPU の構成

4.2 デバイスの構成

CUDA を用いたときの GPU の構成を図 6 に示す。CUDA を利用すると GPU は **grid**, **block**, **thread** の 3 つの単位で管理される。複数の **thread** の集まりを **block** といい、複数の **block** の集まりを **grid** と定義する。1 つの **grid** は 1 つの GPU に対応する。GPU は複数の Streaming Multi Processor (以降 SM) からなり、1 つの SM は 1 つの **block** に対応する。SM は複数の Streaming Processor (以降 SP) からなり、1 つの SP は 1 つの **thread** に対応している。各 **thread** に割り当てられた処理は同時に実行され、並列処理が実現される。

CUDA では実行の際に Warp という単位で同じ命令が行われ同時に処理される。Warp は 32 個の **thread** からなるため、**thread** 数は 32 の倍数の場合が望ましい。Warp 内の各 **thread** が分岐命令でそれぞれが違う方向に分岐を繰り返していくと、並列処理が行われないため GPU の性能低下を招く。このことを Warp divergent といい、GPU で処理を行う際には注意しなければならない。

5. 関連研究

5.1 配列アラインメント

配列アラインメントには、2 つの配列を比較するペアワイズ配列アラインメント (pair-wise sequence alignment) と、3 つ以上の配列を比較する多重配列アラインメント (multiple sequence alignment) がある。また、比較サイズの大きさの面での分類では、配列の全体が

含まれるように比較を行う大域的配列アラインメント (global sequence alignment) と配列の一部で比較を行う局所的配列アラインメント (local sequence alignment) がある。配列アラインメントのうち SW 法は精度を重視するが、精度を犠牲にして速さを重視する手法には BLAST (Basic Local Alignment Search Tool)⁷⁾ や FASTA (FAST-ALL)³⁾ がある。

多重配列アラインメントを求める方法には、ペアワイズ配列アラインメントで使用している DP 法を 3 つの配列に拡張する手法がある⁸⁾。DP 法の拡張を行う手法では、3 配列以上の比較において、必要な計算ステップ数やメモリ量は解析する配列の数に対して、指数関数的に増加するため、文字列長が短い配列しか扱えないという問題がある。よって、一般的に多重配列アラインメントの計算には近似法が用いられている。

多重配列アラインメントのほとんどは近似法を用いている⁹⁾¹⁰⁾¹¹⁾¹²⁾。近似法の多重配列アラインメントは近似法を用いていない手法に比べて、配列アラインメントの精度が劣る。また、大域的配列アラインメントであるため、類似度が最大のもの以外の類似部分を抽出することができない。それに対して、SW 法は局所的配列アラインメントであるため、 n 次元に拡張しても類似度が最大のもの以外の多重配列アラインメントを行うことができる。

5.2 配列アラインメントの GPU 実装

配列アラインメントを GPU で高速化する研究¹³⁾¹⁴⁾¹⁵⁾¹⁶⁾ も多く行われている。SW 法の GPU 実装の多くは、SW 法の文字列テーブルを全て保持しておらず、小さなブロックに分割して保持している場合が多い。これは、SW 法の文字列テーブルはデータ量が膨大になるため、デバイス側からホスト側に文字列テーブルを戻すときのデータ転送に時間がかかるためである。また、多くの実装ではデバイス側で文字列テーブルの計算を実行して、ホスト側でトレースバックを行っている。

5.3 並列計算機向けの SW 法

並列計算機に特化した SW 法には Smith-Waterman-Gotoh (SWG)¹⁷⁾ や Striped Smith-Waterman (SSW)¹⁸⁾ などがある。SWG は、FPGA や GPU などの CPU に比べて複雑な命令処理が苦手なアーキテクチャ向けのアルゴリズムである。一方、SSW は、CPU や Cell.B.E. などの複雑な命令処理が比較的得意であるアーキテクチャ向けのアルゴリズムである。

6. 実 装

6.1 概 要

fNIRS の脳血流量のデータを 3 次元 SW 法を用いて解析を行う。SW 法は文字列の比較

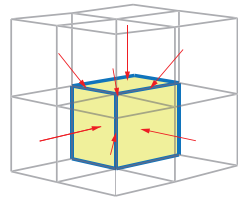


図7 各要素のデータの依存性

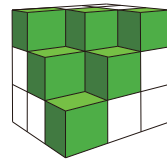


図8 依存性のない要素

のみを行うので、fNIRSの脳血流量データを文字列に変換し、SW法を利用できるようにする。SW法の文字列テーブルをそのまま計算すると時間がかかるので、CUDAを用いてGPUで3次元SW法を実行する。なお、本研究報告で用いた3次元SW法はトレースバックの実装を行っていない。

6.2 時系列データから文字列データへの変換

fNIRSの出力データのような時系列データを、文字列データに変換する手法としては、SAX (Symbolic Aggregation approxXimation)¹⁹⁾ や等間隔領域分割などがある。時系列データを文字列に変換することで、相同性検索を行うことが可能となる。本研究報告では、SAXを用いてfNIRSの脳血流データを文字列に変換し、時系列データの類似部分抽出が行えるようにする。

6.3 GPU実装の概要

3次元SW法を下記に示す方法で実装し、それぞれを比較し評価する。ここでは、3次元SW法の文字列テーブルを4文字毎に区切ったテーブルを、サブテーブルと定義する。

- サブテーブル使用 GPU版3次元SW法 (GPU)
- サブテーブル不使用 GPU版3次元SW法 (GPU)
- サブテーブル不使用 CPU版3次元SW法 (CPU)

図7に3次元SW法の文字列テーブルにおける各要素のデータの依存性を示す。また、図8に3次元SW法の文字列テーブルの斜めの面において依存性のない要素を示す。3次元SW法の文字列テーブルにおいて、それぞれの要素は図7に示すように7近傍の要素とデータが依存関係にある。そのため、図8に示すように斜めの面において、各要素は依存性を持たず、並列に計算が可能である。また、図8に各面ごとに計算可能な要素を示す。このようなSW法の並列性を利用してCUDAで3次元SW法のGPU実装を行う。

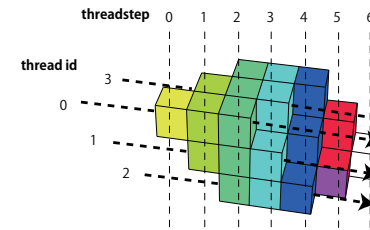


図9 threadの計算手順

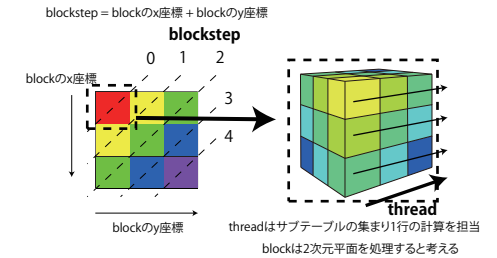


図10 blockの計算手順

6.4 GPUでの実装の手順

6.4.1 文字列テーブルのthread, blockへの割り当て

CUDAにおけるthreadを図9のように割り当てる。一つのthreadで文字列テーブルの一边を全て計算する。図9では、step=3のときに最大7つのスレッドが並列に実行されている。また、一つのthreadで、3次元SW法の文字列テーブルの一边を全て計算するため、blockは、文字列テーブルの次元数が減少し、平面上にblockを割り当てることと同等になる。図10に示すように、blockのx座標とy座標を足したものが現在のstep数と等しいときに、そのブロックが実行される。

6.4.2 サブテーブルの作成

文字列テーブルの計算を高速化するために、図11に示すように、文字列テーブルの要素を4×4×4の64要素毎に区切ったサブテーブルを作成し、割り当てたthreadでそれぞれ処理していく。

通常、デバイス側でメモリを確保するとglobal memoryが確保される、global memoryで文字列テーブルをすべて計算すると、メモリへのアクセスに時間がかかるため、アクセス時間が比較的高速であるshared memoryを用いる。一方、shared memoryはメモリ容量がglobal memoryより小さく、格納するデータの大きさが制限される。よって、global memoryよりアクセスの高速なshared memoryで、64要素という比較的容量の小さいサブテーブルを用いて部分的に文字列テーブルを計算する。GPUは条件分岐処理において処理が低下するので、64要素の値の計算をすべて論理演算を使用し条件分岐を削減した。

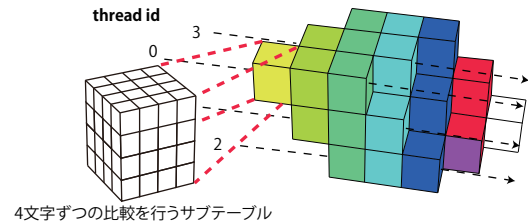


図 11 サブテーブルの作成

7. 評価

7.1 概要

3次元SW法をfNIRSの脳血流量データに基づいて、GPUとCPUで実行時間について性能評価を行う。fNIRSのデータから生成した3つの文字列を、3次元SW法にかけて、文字列テーブルを計算し、CUPS (Cell Updates Per Second:一秒間に文字列テーブルの要素がいくつ計算できるか)について比較検討を行う。SW法の性能評価の指標としてCUPSが一般的に用いられている。比較される文字列の長さを $|A|$ 、比較する文字列の長さの合計を $|B|$ 、実行時間を t 秒とすると、GCUPS (Giga CUPS) は、 $(|A| \times |B|) \div (t \times 10^9)$ で表される。fNIRSのデータの個数、すなわち文字列長が8, 16, 128, 256, 512のときについてそれぞれSW法で文字列テーブルを計算する。本研究報告では文字列テーブルを計算し、ホスト側に文字列テーブルを戻してくるまでの時間を実行時間とする。

3次元SW法を表1の構成のマシンを用いて評価を行う。

7.2 GPU実装の評価

GPUで文字列テーブルを分割して共有メモリに格納し計算を行った場合(GPUI)、GPUでサブテーブルを使用しない場合(GPUN)、CPUでサブテーブルを使用しない場合(CPU)の3次元SW法について評価を行った。結果を図12に示す。

表 1 使用マシン構成

CPU	Intel Core i5-2400 3.10GHz
GPU	GeForce GTX 460
Memory	8GB
OS	Ubuntu 11.04
開発環境	CUDA 4.0
コンパイルオプション	-O3

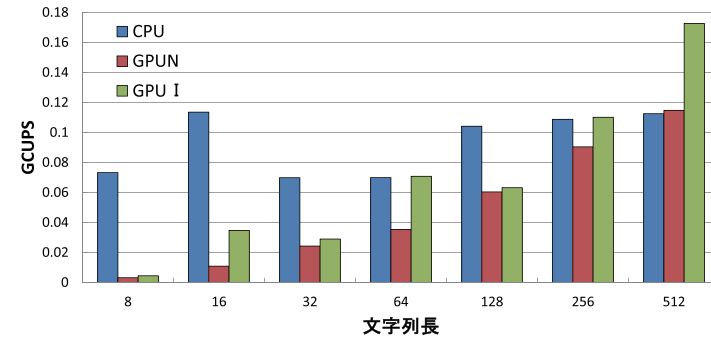


図 12 3次元SW法の評価結果

評価の結果、GPUIではCPUに比べて最大1.53倍、GPUNではCPUに比べて最大1.01倍、GPUIはGPUNに比べて最大1.50倍の性能向上がみられた。また、GPUIは最大で0.172GCUPSを示した。また最低で0.004GCUPSを示した。

図12より、文字列長が短い程CPUの方が高速な処理を行っていることがわかる。これは、文字列テーブルの計算時間より、ホスト-デバイス間のデータの転送に時間を要することによる。また、文字列長が長くなるにつれ、CPUに比べてGPU実装の方が処理性能が高くなることわかる。

GPUNに比べてGPUIの方が文字列長を長くした場合の、処理性能の向上率が大きいこともわかる。これは、文字列長が増えるほど文字列テーブルへのメモリアクセスが増加するという性質をもつSW法を実行ときのアルゴリズムの違いから発生することである。GPUIはサブテーブルに分割して一度に複数のデータを、メモリアクセスの速いshared memoryで計算し、値の格納をまとめて行うのに対して、GPUNでは文字列テーブルを計算する度、メモリアクセスの遅いglobal memoryからデータを取得しているためである。

次に、GPUNについてGPUの枚数と性能の関係を図13に示す。2枚のGPUについてそれぞれ同じ文字列の比較を行っている。評価の結果GPUを2枚にしたところ、GPUI枚のときに比べて約1.8倍の性能向上がみられた。図13より、文字列長によって性能は変わらないことがわかった。また、単純に処理性能が2倍とならないのは、転送にかかるオーバーヘッドによるものと考えられる。

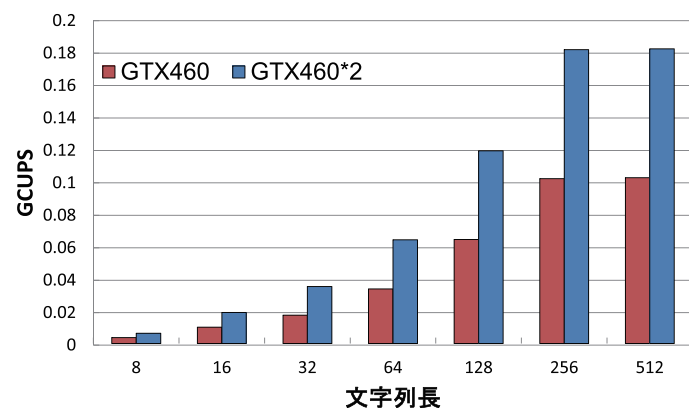


図 13 GPUN における GPU の枚数による 3 次元 SW 法の評価結果

8. まとめと今後の展開

本研究報告では、3次元 SW 法の提案を行い、fNIRS の脳血流データに基づいて評価を行った。結果として工夫ありの GPU 実装は工夫なしの CPU の実装に比べて最大 1.53 倍の性能向上がみられた。また、2 枚の GPU を用いると文字列長の長さに関係なく約 1.8 倍の性能向上がみられた。

今後の展開としては並列計算機向けの SW 法の導入による処理の高速化、また、さらに GPU の枚数を増やして性能の変化を検討することが考えられる。

参考文献

- 1) Owens, J.D., Luebke, D., Govindaraju, N., Harris, M., Kruger, J., Lefohn, A. and Purcell, T.J.: A Survey of General-Purpose Computation on Graphics Hardware, *Eurographics 2005, State of the Art Reports*, pp.21–51 (2005).
- 2) Smith, T. and Waterman, M.: Identification of Common Molecular Subsequences, *Journal of Molecular Biology*, Vol.147, No.1, pp.195–197 (1981).
- 3) Mount, D.W.: バイオインフォマティクス ゲノム配列から機能解析へ第 2 版, メディカル・サイエンス・インターナショナル (2005).
- 4) 松下晋, 中川匡弘: 光トポグラフィーによる感性情報解析 (ヒューマンコミュニケー

- ション), 電子情報通信学会論文誌. A, 基礎・境界, Vol.88, No.8, pp.994–1001 (2005-08-01).
- 5) 廣安知之, 西井琢真, 吉見真聡, 三木光範, 横内久猛: 相同性検索を用いた 2 つの時系列データからの類似部分抽出手法と DTW による類似部分の評価, 情報処理学会研究報告. MPS, 数理モデル化と問題解決研究報告, Vol.2010, No.24, pp.1–6 (2010-09-21).
- 6) Henikoff, S. and Henikoff, J.: Amino acid substitution matrices from protein blocks., *Proceedings of the National Academy of Sciences of the United States of America (PNAS)*, Vol.89, pp.10915–10919 (1992).
- 7) NCBI: BLAST Basic Local Alignment Search Tool (2009).
- 8) Liu, Y., Schmidt, B. and Maskell, D.L.: MSA-CUDA: Multiple Sequence Alignment on Graphics Processing Units with CUDA, *Application-Specific Systems, Architectures and Processors, IEEE International Conference on*, Vol.0, pp.121–128 (online), DOI:http://doi.ieeecomputersociety.org/10.1109/ASAP.2009.14 (2009).
- 9) Thompson, J.D., Higgins, D.G. and Gibson, T.J.: CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice, *Nucleic Acids Research*, Vol.22, No.22, pp.4673–4680 (1994).
- 10) Katoh, K., Misawa, K., Kuma, K. and Miyata, T.: MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform, *Nucleic Acids Research*, Vol.30, No.14, pp.3059–3066 (2002).
- 11) Bray, N. and Pachter, L.: MAVID: constrained ancestral alignment of multiple sequences, *Genome Research*, Vol.14, No.4, pp.693–699 (2004).
- 12) Notredame, C., Higgins, D.G. and Heringa, J.: T-Coffee: A Novel Method for Fast and Accurate Multiple Sequence Alignment, *Journal of Molecular Biology*, Vol.14, No.4, pp.693–699 (1994).
- 13) 土肥慶亮, 柿本雄, 柴田裕一郎, 濱田剛, 小栗清: GPU クラスタにおける Smith-Waterman アルゴリズムの実装手法の提案, *SACIS2010*, Vol.2010, No.5, pp.209–216 (2010).
- 14) 宗川裕馬, 伊野文彦, 荻原兼一: 統合開発環境 CUDA を用いた GPU での配列ライメントの高速化手法, 情報処理学会研究報告, Vol.114, No.19, pp.13–18 (2008).
- 15) de OSandes, E.F. and de Melo, A. C. M.A.: Smith-Waterman Alignment of Huge Sequences with GPU in Linear Space, *2011 IEEE International Parallel & Distributed Processing Symposium*, Vol.25, pp.1199–1211 (2011).
- 16) Liu, Y., Schmidt, B. and Maskell, D.L.: CUDASW++2.0: enhanced Smith-Waterman protein database search on CUDA-enabled GPUs based on SIMT and virtualized SIMD abstractions, *BMC Research Notes*, Vol.3, No.1, pp.1–12 (2010).
- 17) Gotoh, O.: An Improved Algorithm for Matching Biological Sequences, *Journal of Molecular Biology*, Vol.162, pp.705–708 (1981).

- 18) Farrar, M.: Striped Smith–Waterman speeds database searches six times over other SIMD implementations, *Bioinformatics*, Vol.23, pp.156–161 (2007).
- 19) Lin, J., Keogh, E., Lonardi, S. and Chiu, B.: A Symbolic Representation of Time Series, with Implications for Streaming Algorithms, *In Proceedings of the 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery*, ACM Press, pp.2–11 (2003).